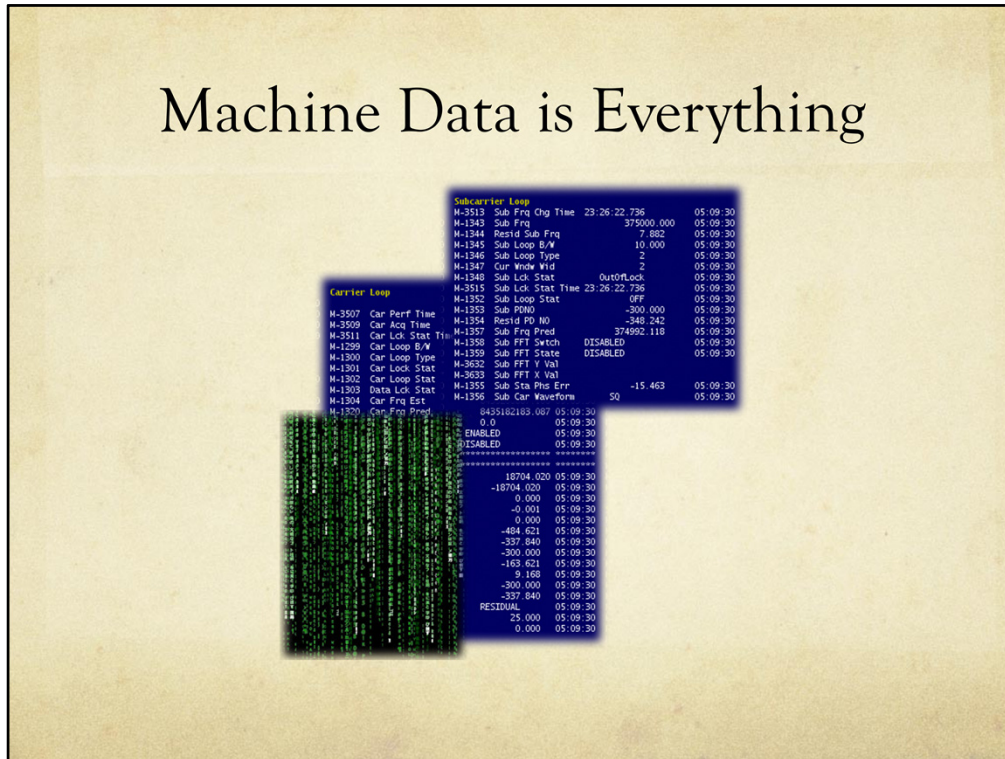Hi everybody. Today I'm going to talk about EPOXI's experience with a Big Data tool, Splunk. The goal of this experiment was to try and improve our team's operational intelligence, or in another words, the ability for us to see trends and high level information in the machine data

- The reason I started down this path to begin with was because I'm quite new to mission operations. I joined the team about 2 years ago having worked mostly in web and software development, with an interest in data visualization.
- So when I began learning to fly the spacecraft, I was instantly overwhelmed by the sheer amount of data I had to filter through. In this image, you can see one of our telemetry displays. Unlike my very seasoned co-workers, I had a difficult time getting an overall understanding of what was going on just by glancing at these pages. If a problem occurred, I wouldn't notice until later because the important information would be hidden in pages and pages of data.
- . Once I became more familiar with the telemetry data and displays, I was still pretty inefficient at reading it quickly because I couldn't customize the interface to my liking, and the data was more low-level than I needed to see.
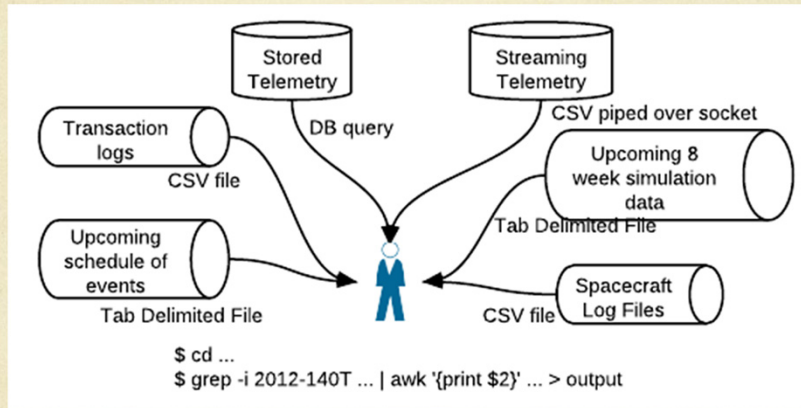
I had trouble with other things beyond the telemetry system as well. when analyzing historic spacecraft data, I always wasted a lot of time collecting data. On top of that, I would spend a non-trivial amount of time scripting one-time-use tools to do a quick analysis. For example, there were times when I wanted to make sure our file downlink rate was what I expected. This is important in order to know how many files we can downlink in a given time. In order to do this, I would have to compare our prediction data with the timestamps of files. This picture shows the two formats of the information – a flat file with a bunch of text, and a long directory listing of when the files were actually received on the ground. Doing this analysis was a little annoying.

So, to put the problem in concrete terms, the trouble I had stemmed from the fact that our data was scattered all over the place in difficult to read formats. Beyond real time telemetry data, we have transaction logs, downlinked images, spacecraft logs, schedules, and more. These were stored on a number of different servers and databases, and all had varying formats.

And, as I said before, parsing this data always involved some serious command line munging.

# Splunk?



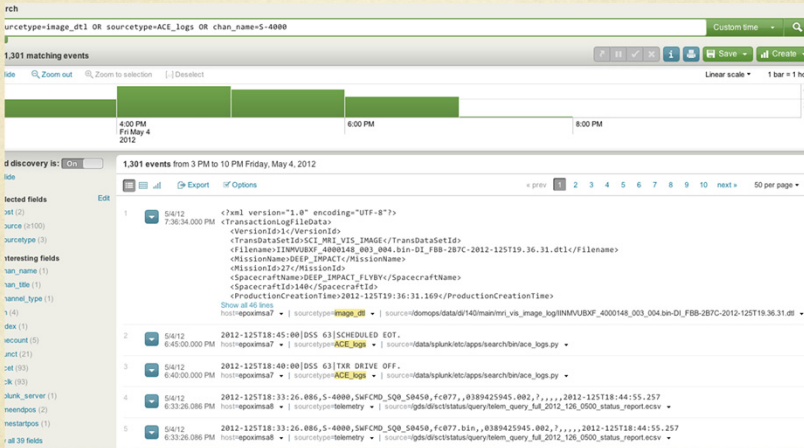So what is Splunk? Instead of giving the technical details, which you can find online, I'll tell you what it did for me.  Splunk slapped everything into one place, with one uniform format, and gave me the ability to forget about all these annoying details of where it is, how to parse it, and all that. Instead, I only need to interact with Splunk to find the data I need. This sounds simple and obvious, but it's surprising what you can do once you all of your data is indexed in one place.

By having your data organized, querying becomes much easier. Let's say that I want to search telemetry for a sensor_name "temp_1" and to return all data that is at most five minutes old. And because Splunk can hook into a real-time stream, this data will always be up-to-date.

# Data Mashup



Extending the previous example, I can now aggregate all types of data into one view based in time. In this picture, I've got transaction logs, telemetry, and downlinked files all in one page, organized by time. Even though the raw data looks completely than this, I've defined interfaces that transform it into this uniform format. This gives me a more complete picture for the question "what was the spacecraft doing at this particular time?" And because querying data is simple, I can start with a big block of data and whiddle it down to what I need, rather than hunting around for the individual pieces of data that I need.
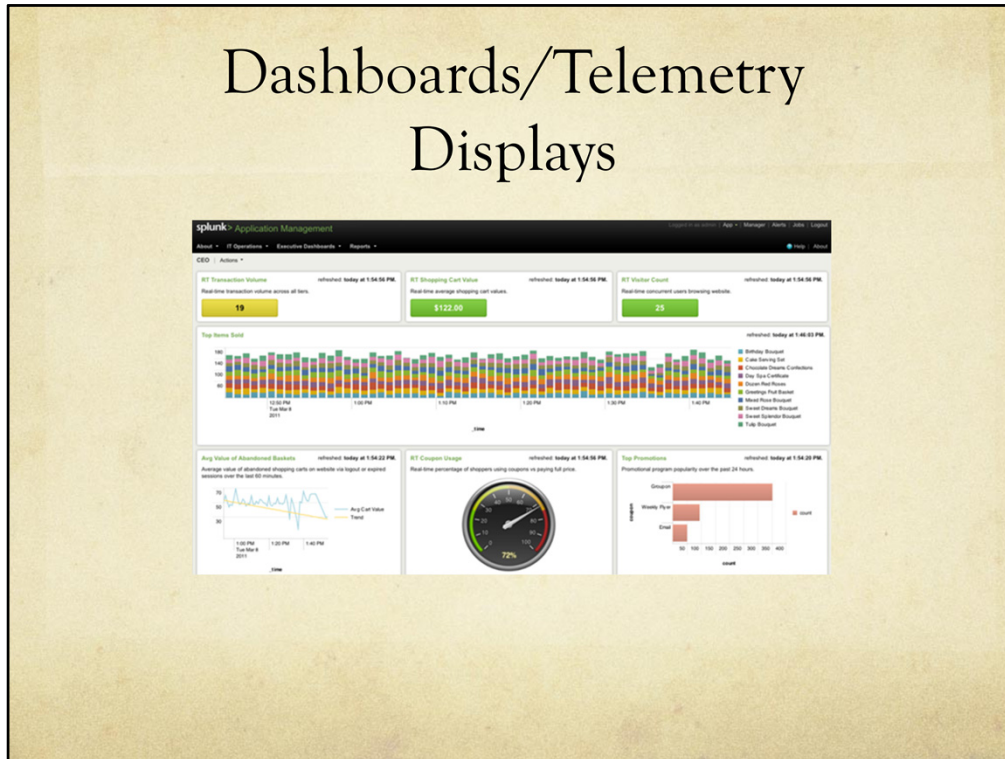
# What can I do now?

- Unix-like queries
  - Sensor_name=temp_1 | head 5
  - Sensor_name=temp_1 | sort –temp| head 5
  - Sensor_name=temp_1 | delta temp as tmp_chg | where tmp_chg > 10

When we have all the data we need, we can begin widdling down the data with Splunk's Unix-like search syntax. These three examples highlights my trial-and-error attempts to find large temperature changes. I begin by showing the first 5 temperatures, only to find that they're sorted chronologically, rather than from highest temperatures to lowest temperatures. The next line shows sorting temperatures by their values, but I find that that's not really what I want either. I want to know the delta temperatures between readings.  Looking through Splunk's user manual, I find the delta function, which lets me dynamically generate new information to use in my query. With that extra piece of information, I can now return only the telemetry readings where the temperature changed by at least 10.

One other useful feature I'll mention is that all of these queries can be run through Splunk's API. So any scripting language you can think of can plug right in and make these queries. This gives us the ability to build a lot of new tools.

Dashboards/Telemetry Displays

With this query language, you can start building some pretty cool tools.  If you look carefully, you'll notice this telemetry display, or dashboard, has absolutely nothing to do with spacecrafts.  However, it's a good example of the user interface and visualizations Splunk can provide.  This in-browser dashboard hooks right into Splunk's search system, and you can begin building some impressive displays with only a few clicks of a mouse.

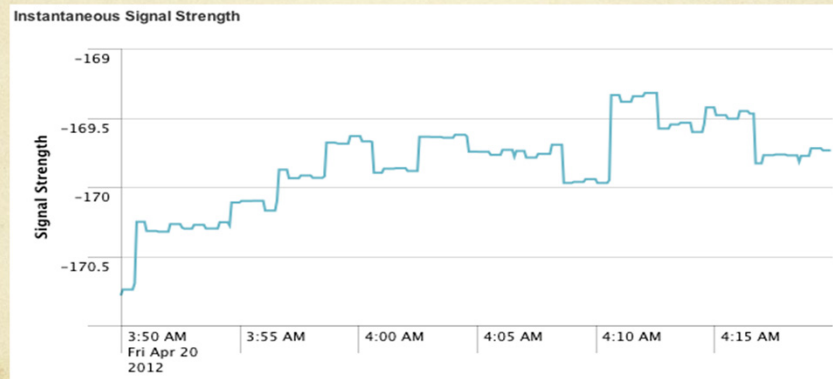Now I'll show you some of my dashboards, which are not nearly as pretty as this.

This image shows one dashboard element I made with only one search query. As I mentioned before, I often check to see how long it takes to downlink an image. Here, as the images come down in real-time, Splunk's system continuously refreshes a table containing each file and their downlink time. Again, this was formed with just a few search queries similar to those on the previous slide.
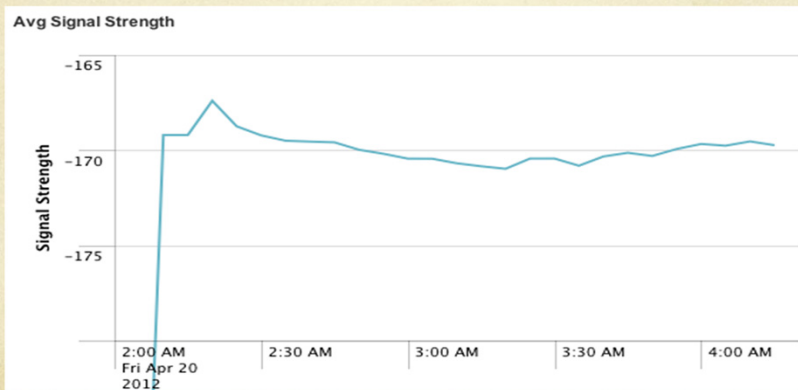
# Quick Plotting

○ `Channel_name=signal_strength earliest=-40m | timechart span=5s first(digital_value)`



Splunk has a powerful graphing module that allows you to set up real-time plots very quickly. Here, we're plotting signal strength over time by specifying the period that we'd like to graph over (last 40 minutes) and taking one value every 5 seconds. The problem with this graph is that we cannot see the trend – the timespan is too short, and the data is too noisy. In our old system, this would take a bit of scripting to fix.
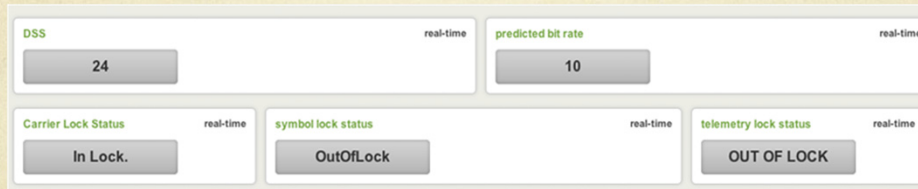
## Quick Plotting

```
sourcetype=DSN telemetry channel name=signal strength earliest=-2h |
timechart span=10m median(digital value)
```
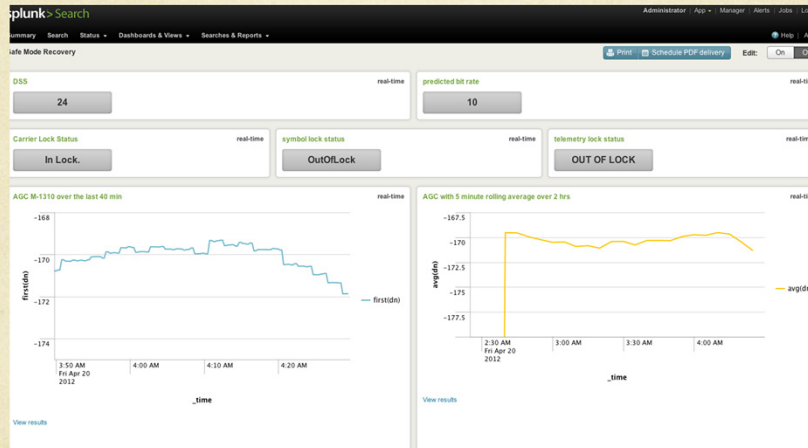
Avg Signal Strength

However, all it takes with Splunk is to alter the search syntax to span the past two hours, bucket the data into 10 minute chunks and take the median.  Now we see the trend much more clearly.

Along with the graphing modules, Splunk also has intuitive user interfaces for building "dashboards", which are equivalent to telemetry displays. While this telemetry display isn't novel by any means, it was built in less than five minutes with a simple point-and-click interface. Because Splunk also has LDAP authentication built in, each user can build their own unique telemetry display or pick from a global set of displays. Because these are all in-browser displays, they can be pulled up virtually anywhere with any system.

So here is a complete picture of one of my dashboards, with the elements that I described earlier.  Not nearly as pretty as the promotional screenshot, but considering this took me 10 minutes to make, I'm pretty pleased.

As a side note, one useful feature that our current system lacks is back-fill. Meaning, when I pull up this page, it will populate graphs immediately without having to wait for more data to stream it.

# Alerting System

- Query until your expected result appears…if not, send an alert
  - Expect data from spacecraft by beginning-of-track + 30 minutes. If not, send alert
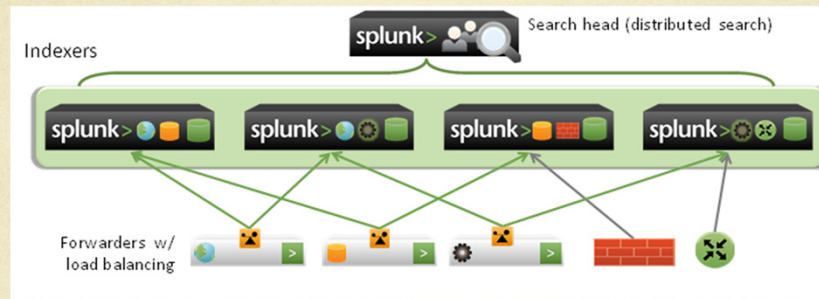  - $ search telemetry earliest=-5m | head 1



Because searches can be run continuously, Splunk makes it easy to set up an alerting system in case your data isn't showing what you expect it to. One simple example is sending an alert if the spacecraft isn't talking to us when we expect it to. To set this up, you would activate the alert based on some tracking schedule. Then, you specify the search to begin at the beginning of the track, and to trigger if no results have appeared

Distributed Architecture and Scaling

Now, you might be wondering what Splunk actually looks like. I won't get too much into detail, but to put it into real terms, Splunk has an architecture that allows it to be distributed across as many machines as you need. Obviously, if you're dealing with terabytes of data a day, one machine won't cut it. By designating many machines to only index data or to only process search queries, you can keep your data behind firewalls or on specific machines without worrying about data privacy issues while simultaneously increasing the throughput of your system.

# Conclusion: What did Splunk do for us?

- Organized our data

- Provided a new platform to develop tools and manipulate data

- Gave us new user interfaces to our real-time data